# CR-SAM: Curvature Regularized Sharpness-Aware Minimization

Tao Wu[1]; Tie Luo[1*]; Donal C. Wunsch II[2]
Missouri University of Science and Technology
[1]Department of Computer Science, [2]Department of Electrical and Computer Engineering

## BACKGROUND

- **Sharpness-Aware Minimization (SAM)** directs the search for model parameters within flat regions by minimizing the maximum loss in the vicinity of weight $\boldsymbol{w}$ within a radius $\rho$

$$\min_{\boldsymbol{w}} L^{\mathrm{SAM}}(\boldsymbol{w}) \text{ where } L^{\mathrm{SAM}}(\boldsymbol{w}) = \max_{\|\boldsymbol{v}\|_2 \leq 1} L_{\mathcal{S}}(\boldsymbol{w} + \rho \boldsymbol{v})$$

- SAM employs one-step gradient ascent to approximate the inner maximization

$$\nabla L^{\mathrm{SAM}}(\boldsymbol{w}) \approx \nabla L_{\mathcal{S}}\left(\boldsymbol{w} + \rho \frac{\nabla L_{\mathcal{S}}(\boldsymbol{w})}{\|\nabla L_{\mathcal{S}}(\boldsymbol{w})\|_2}\right)$$

## PITFALLS IDENTIFIED (by us)

- **Deterioration** of one-step sharpness approximation as training proceeds.

  We define *approximation ratio* (AR):

$$\mathrm{AR} = \mathbb{E}_{(x,y) \sim D}\left[\frac{\ell\left(f(x; \boldsymbol{w} + \boldsymbol{\delta}), y\right) - \ell\left(f(x; \boldsymbol{w}), y\right)}{\ell\left(f(x; \boldsymbol{w} + \boldsymbol{\delta^*}), y\right) - \ell\left(f(x; \boldsymbol{w}), y\right)}\right]$$

  where $\boldsymbol{\delta}$ represents *one-step* gradient ascent perturbation, and $\boldsymbol{\delta^*}$ denotes the optimal perturbation which is calculated with 20-step gradient ascent.

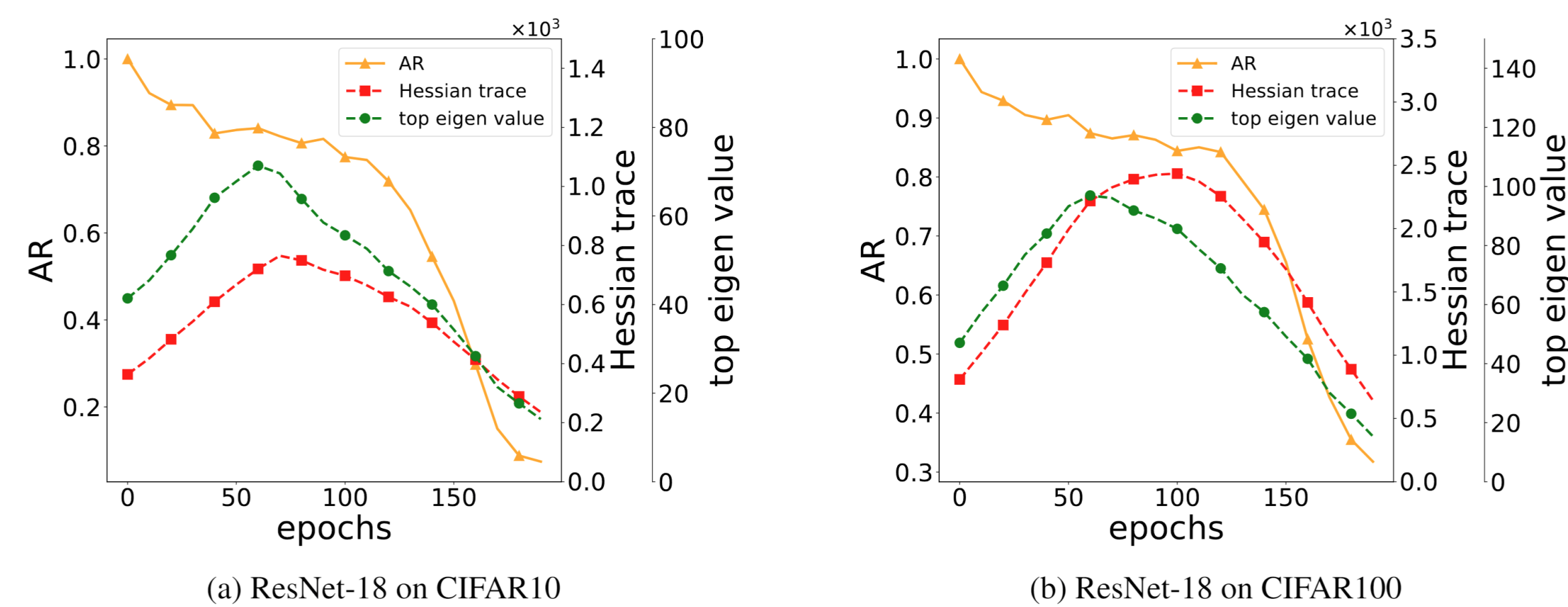- **Conventional Hessian fails to capture true loss landscape curvature**



Figure 1. Evolution of approximation ratio (AR), Hessian trace, and top eigenvalue of Hessian during SAM training on CIFAR10 and CIFAR100 datasets. The continuously decreasing AR indicates an enlarging curvature whereas both of the Hessian-based curvature metrics, Hessian trace and the top eigenvalue (which are expected to continuously increase) fail to capture the true curvature of model loss landscape.

## RESOURCES AND CONTACT

- **Paper**: https://arxiv.org/abs/2312.13555
- **Code**: https://github.com/TrustAIoT/CR-SAM
- **Contact**: wuta@mst.edu, tluo@mst.edu, dwunsch@mst.edu

## METHODS

- **Propose a new metric for accurate curvature characterization, called *normalized Hessian trace*:**

$$\mathcal{C}(\boldsymbol{w}) = \frac{\mathrm{Tr}\left(\nabla^2 L_{\mathcal{S}}(\boldsymbol{w})\right)}{\|\nabla L_{\mathcal{S}}(\boldsymbol{w})\|_2}$$

- **Characterizes curvature faithfully**
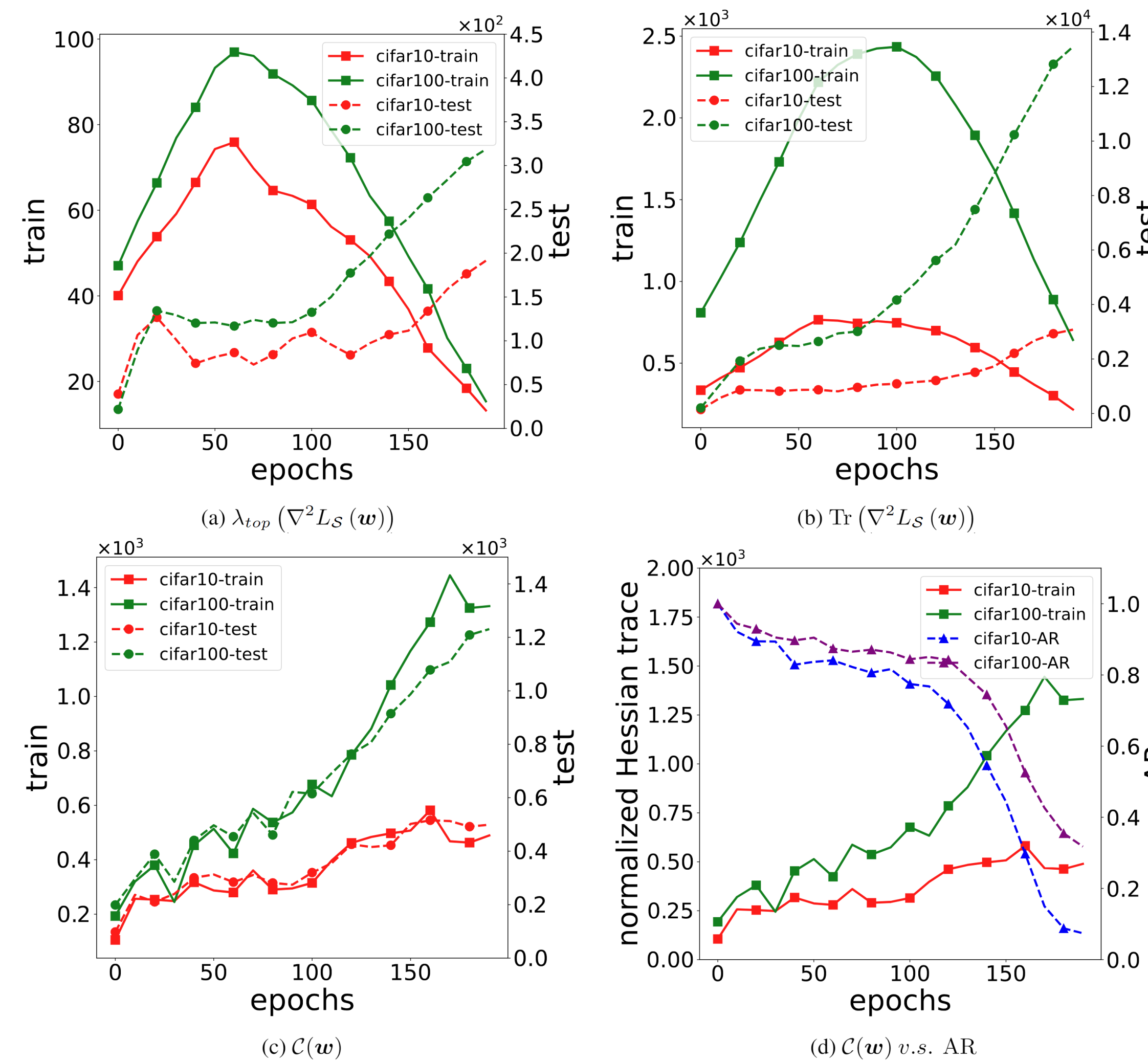- **Consistent trends on both training and test sets**



Figure 2. The **normalized Hessian trace (d)** indicates a growing curvature throughout training which explains the phenomenon of decreasing AR (approximation ratio). In addition, it behaves consistently on both train and test sets (**c**; as opposed to a & b which are conventional Hessian).

- **Curvature Regularized Sharpness-Aware Minimization (CR-SAM)**

$$R_c(\boldsymbol{w}) = \alpha \log \mathrm{Tr}\left(\nabla^2 L_{\mathcal{S}}(\boldsymbol{w})\right) + \beta \log \|\nabla L_{\mathcal{S}}(\boldsymbol{w})\|_2$$

$$\min_{\boldsymbol{w}} L^{\mathrm{CR\text{-}SAM}}(\boldsymbol{w})$$

where $L^{\mathrm{CR\text{-}SAM}}(\boldsymbol{w}) = L^{\mathrm{SAM}}(\boldsymbol{w}) + R_c(\boldsymbol{w})$

Solving it with *finite difference method*:

$$R_c(\boldsymbol{w}) = \mathbb{E}_{\boldsymbol{v} \sim N(0,I)}\Big[\alpha \log\left(L_{\mathcal{S}}(\boldsymbol{w} + \rho\boldsymbol{v}) + L_{\mathcal{S}}(\boldsymbol{w} - \rho\boldsymbol{v})\right.$$
$$\left. - 2L_{\mathcal{S}}(\boldsymbol{w})\right) + \beta \log\left(L_{\mathcal{S}}(\boldsymbol{w} + \rho\boldsymbol{v}) - L_{\mathcal{S}}(\boldsymbol{w} - \rho\boldsymbol{v})\right)\Big]$$

## RESULTS

|  |  | CIFAR-10 | | | CIFAR-100 | | |
|---|---|---|---|---|---|---|---|
| Model | Aug | SGD | SAM | CR-SAM | SGD | SAM | CR-SAM |
| ResNet-18 | Basic | 95.29±0.16 | 96.46±0.18 | **96.95**±0.13 | 78.34±0.22 | 79.81±0.18 | **80.76**±0.21 |
|  | Cutout | 95.96±0.13 | 96.55±0.15 | **97.01**±0.21 | 79.23±0.13 | 80.15±0.17 | **81.26**±0.19 |
|  | AA | 96.33±0.15 | 96.75±0.18 | **97.27**±0.12 | 79.05±0.17 | 81.26±0.21 | **82.11**±0.22 |
| ResNet-101 | Basic | 96.35±0.12 | 96.51±0.16 | **97.14**±0.11 | 80.54±0.13 | 82.11±0.12 | **83.03**±0.17 |
|  | Cutout | 96.56±0.18 | 96.95±0.13 | **97.51**±0.24 | 81.26±0.21 | 82.39±0.27 | **83.46**±0.16 |
|  | AA | 96.78±0.14 | 97.11±0.16 | **97.76**±0.16 | 81.83±0.37 | 83.25±0.47 | **84.19**±0.23 |
| WRN-28-10 | Basic | 95.89±0.21 | 96.81±0.26 | **97.36**±0.15 | 81.84±0.15 | 83.15±0.14 | **84.45**±0.09 |
|  | Cutout | 96.89±0.07 | 97.55±0.16 | **97.98**±0.21 | 81.96±0.40 | 83.47±0.15 | **84.48**±0.13 |
|  | AA | 96.93±0.12 | 97.59±0.06 | **97.94**±0.08 | 82.16±0.11 | 83.69±0.26 | **84.74**±0.21 |
| PyramidNet-110 | Basic | 96.27±0.13 | 97.34±0.13 | **97.89**±0.08 | 83.27±0.12 | 84.89±0.09 | **85.68**±0.14 |
|  | Cutout | 96.79±0.13 | 97.61±0.21 | **98.08**±0.11 | 83.43±0.26 | 84.97±0.17 | **85.86**±0.21 |
|  | AA | 96.97±0.08 | 97.81±0.13 | **98.26**±0.11 | 84.59±0.08 | 85.76±0.23 | **86.58**±0.14 |

**Table 1.** Classification accuracy on CIFAR-10 and CIFAR-100 datasets.

| **Optimizer** | $\|\nabla L_{\mathcal{S}}(\boldsymbol{w})\|_2$ | $\mathrm{Tr}\left(\nabla^2 L_{\mathcal{S}}(\boldsymbol{w})\right)$ | $\mathcal{C}(\boldsymbol{w})$ |
|---|---|---|---|
| SGD | 19.97 ±0.52 | 32673 ±1497 | 1674 ±78 |
| SAM | 11.51 ±0.31 | 14176 ±327 | 1193 ±59 |
| CR-SAM | **8.26** ±0.19 | **7968** ±145 | **884** ±23 |

**Table 2.** Model geometry (characterized by 3 metrics) of ResNet-18 trained with SGD, SAM and CR-SAM. Values are computed on test set. It shows the CR-SAM optimizer achieves the minimal in all cases.
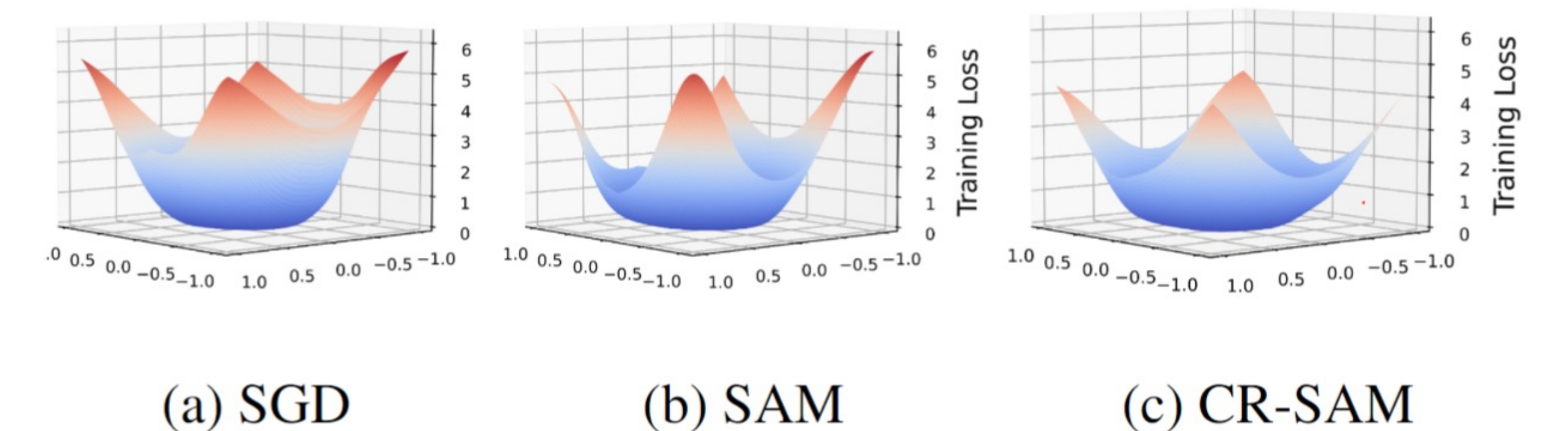


(a) SGD          (b) SAM          (c) CR-SAM

**Figure 3.** CR-SAM yields flatter loss landscape.



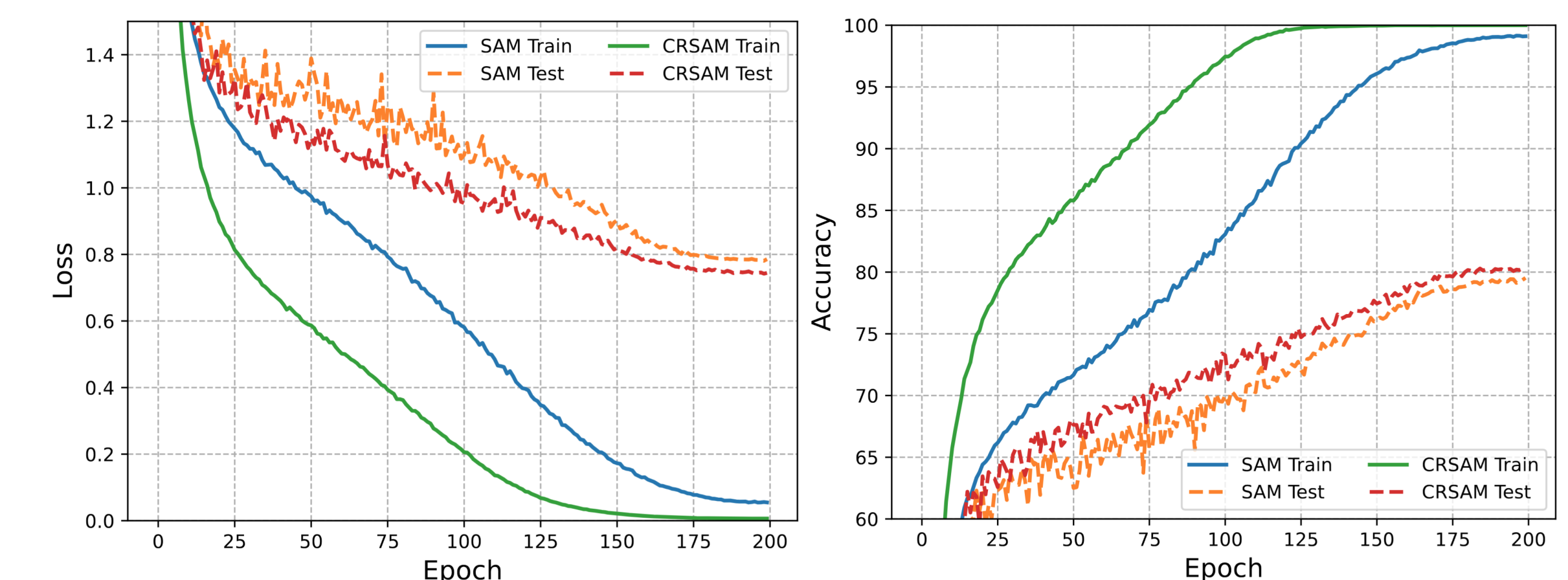**Figure 4.** CR-SAM achieves much faster and stabler convergence. This can be explained by the fact that CR-SAM discourages excessive curvature and thus reduces optimization complexity, making local minima easier to reach.